

STATISTICS FOR
DATA SCIENCE

Top 50 Interview Questions & Answers



WITH EXPLANATORY
YOUTUBE LINKS

1) What is Inferential Statistics?

Inferential statistics is a statistical method that deduces from a small but representative sample the characteristics of a bigger population. In other words, it allows the researcher to make assumptions about a wider group, using a smaller portion of that group as a guideline.

2) What is the difference between Population and Sample?

From the population we take a sample. We cannot work on the population either due to computational costs or due to availability of all data points for the population. From the sample we calculate the statistics. From the sample statistics we conclude about the population.

3) What is the relationship between mean and median in normal distribution?

In the normal distribution mean is equal to median

4) What is an outlier?

An outlier is an abnormal value (It is at an abnormal distance from the rest of the data points).

5) What can I do with Outliers?

- Remove outlier

When we know the data-point is wrong (negative age of a person)

When we have lots of data

We should provide two analyses. One with outliers and another without.

- Keep outlier

When there are lot of outliers (skewed data)

When results are critical

When outliers have meaning (fraud data)

[\(click on the youtube button to watch explanatory video\)](#)

6) What is the difference between population parameters and sample statistics?

- Population parameters are:

Mean = μ

Standard deviation = σ

- Sample statistics are:

Mean = \bar{x}

Standard deviation = s

7) What is the difference between inferential statistics and descriptive statistics?

- Descriptive statistics – provides exact and accurate information.
- Inferential statistics – provides information of a sample and we need inferential statistics to reach to a conclusion about the population.

8) What is the difference between population and sample in inferential statistics?

- From the population we take a sample. We cannot work on the population either due to computational costs or due to availability of all data points for the population.
- From the sample we calculate the statistics
- From the sample statistics we conclude about the population

9) What are descriptive statistics?



Descriptive statistic is used to describe the data (data properties)

5-number summary is the most commonly used descriptive statistics

10) Most common characteristics used in descriptive statistics?

- Central Tendency- middle of the data. Mean / Median / Mode are the most commonly used as measures.

Mean – average of all the numbers

Median – the number in the middle

Mode – the number that occurs the most. The disadvantage of using Mode is that there may be more than one mode.

- Spread – How the data is dispersed. Range / IQR / Standard Deviation / Variance are the most commonly used as measures.

Range = Max – Min

Inter Quartile Range (IQR) = $Q3 - Q1$

Standard Deviation (σ) = $\sqrt{(\sum (x-\mu)^2 / n)}$

Variance = σ^2

- Shape – the shape of the data can be symmetric or skewed

Symmetric – the part of the distribution that is on the left side of the median is same as the part of the distribution that is on the right side of the median

Left skewed – the left tail is longer than the right side

Right skewed – the right tail is longer than the left side

- Outlier – An outlier is an abnormal value

Keep the outlier based on judgement

Remove the outlier based on judgement

11) What are 5 Point Summary (or 5 number summary)?



- Low extreme (minimum)
- Lower quartile ($Q1$)
- Median
- Upper quartile ($Q3$)
- Upper extreme (maximum)

12) What are the measures of central tendency? When to use which one?

Measures of central Tendency are – Mean, Median & Mode.

1. Mean is the most frequently used measure of central tendency and generally considered the best measure of it. However, there are some situations where either median or mode are preferred.

2. Median is the preferred measure of central tendency when:

- There are a few extreme scores in the distribution of the data.
(NOTE: Remember that a single outlier can have a great effect on the mean).

- There are some missing or undetermined values in your data.
 - There is an open ended distribution (For example, if you have a data field which measures number of children and your options are 0, 1, 2, 3, 4, 5 or "6 or more," then the "6 or more field" is open ended and makes calculating the mean impossible, since we do not know exact values for this field)
 - You have data measured on an ordinal scale.
3. Mode is the preferred measure when data are measured in a nominal (and even sometimes ordinal) scale.

13) What is CLT? What is its significance?



"The central limit theorem states that the sampling distribution of the sample mean approaches a normal distribution as the sample size gets larger no matter what the shape of the population distribution." The central limit theorem is important because it is used in hypothesis testing and also to calculate confidence intervals.

(Please click for more reference:

<https://towardsdatascience.com/understanding-the-central-limit-theorem-e598158cc5da>)

14) What is p-value?



The P value, or calculated probability, is the probability of finding the observed, or more extreme, results when the null hypothesis (H_0) of a study question is true - the definition of 'extreme' depends on how the hypothesis is being tested. P is also described in terms of rejecting H_0 when it is actually true, however, it is not a direct probability of this state.

When you perform a hypothesis test in statistics, a p-value helps you determine the significance of your results. Hypothesis tests are used to test the validity of a claim that is made about a population. This claim that's on trial, in essence, is called the null hypothesis. The alternative hypothesis is the one you would believe if the null hypothesis is concluded to be untrue.

The evidence in the trial is your data and the statistics that go along with it. All hypothesis tests ultimately use a p-value to weigh the strength of the evidence (what the data are telling you about the population). The p-value is a number between 0 and 1 and interpreted in the following way:

- A small p-value (typically ≤ 0.05) indicates strong evidence against the null hypothesis, so you reject the null hypothesis.
- A large p-value (> 0.05) indicates weak evidence against the null hypothesis, so you fail to reject the null hypothesis.
- p-values very close to the cut-off (0.05) are considered to be marginal (could go either way). Always report the p-value so your readers can draw their own conclusions.

15) What are z and t stats?

Z-test: In a z-test, the sample is assumed to be normally distributed. A z-score is calculated with population parameters such as “population mean” and “population standard deviation” and is used to validate a hypothesis that the sample drawn belongs to the same population.

Null: Sample mean is same as the population mean

Alternate: Sample mean is not same as the population mean

The statistics used for this hypothesis testing is called z-statistic, the score for which is calculated as

$$z = (x - \mu) / (\sigma / \sqrt{n}), \text{ where}$$

x = sample mean

μ = population mean

σ / \sqrt{n} = population standard deviation

If the test statistic is lower than the critical value, accept the hypothesis or else reject the hypothesis

- T-test:

A t-test is used to compare the mean of two given samples. Like a z-test, a t-test also assumes a normal distribution of the sample. A t-test is used when the population parameters (mean and standard deviation) are not known.

There are three versions of t-test

1. Independent samples t-test which compares mean for two groups
2. Paired sample t-test which compares means from the same group at different times

3. One sample t-test which tests the mean of a single group against a known mean. The statistic for this hypothesis testing is called t-statistic, the score for which is calculated as

$$t = (x_1 - x_2) / (\sigma / \sqrt{n_1} + \sigma / \sqrt{n_2}), \text{ where}$$

x_1 = mean of sample 1

x_2 = mean of sample 2

n_1 = size of sample 1

n_2 = size of sample 2

16) What are hypothesis testing?

A statistical hypothesis is an assumption about a population parameter. This assumption may or may not be true. Hypothesis testing refers to the formal procedures used by statisticians to accept or reject statistical hypotheses.

Statistical Hypotheses:

The best way to determine whether a statistical hypothesis is true would be to examine the entire population. Since that is often impractical, researchers typically examine a random sample from the population. If sample data are not consistent with the statistical hypothesis, the hypothesis is rejected.

There are two types of statistical hypotheses.

Null hypothesis:

The null hypothesis, denoted by H_0 , is usually the hypothesis that sample observations result purely from chance.

Alternative hypothesis:

The alternative hypothesis, denoted by H_1 or H_a , is the hypothesis that sample observations are influenced by some non-random cause.

For example, suppose we wanted to determine whether a coin was fair and balanced. A null hypothesis might be that half the flips would result in Heads and half, in Tails. The alternative hypothesis might be that the number of Heads and Tails would be very different. Symbolically, these hypotheses would be expressed as

- $H_0: P = 0.5$
- $H_a: P \neq 0.5$

Suppose we flipped the coin 50 times, resulting in 40 Heads and 10 Tails. Given this result, we would be inclined to reject the null hypothesis. We would conclude, based on the evidence, that the coin was probably not fair and balanced.

Other terms related to Hypothesis Testing are:

- Level of significance: Refers to the degree of significance in which we accept or reject the null-hypothesis. 100% accuracy is not possible for accepting or rejecting a hypothesis, so we therefore select a level of significance that is usually 5%.
- Type I error: When we reject the null hypothesis, although that hypothesis was true. Type I error is denoted by alpha. In hypothesis testing, the normal curve that shows the critical region is called the alpha region.
- Type II errors: When we accept the null hypothesis but it is false. Type II errors are denoted by beta. In Hypothesis testing, the normal curve that shows the acceptance region is called the beta region.
- Power: Usually known as the probability of correctly accepting the null hypothesis. $1 - \beta$ is called power of the analysis.
- One-tailed test: When the given statistical hypothesis is one value like $H_0: \mu_1 = \mu_2$, it is called the one-tailed test.
- Two-tailed test: When the given statistics hypothesis assumes a less than or greater than value, it is called the two-tailed test

17) How do you determine outliers?

The different ways of determining the outliers:

1. Simply sort your data sheet for each variable and then look for unusually high or low values.
2. Graphing your data: Using Boxplots, Histogram and Scatterplots.

Boxplot: We use boxplots when we have groups in our data. Boxplots display asterisks or other symbols on the graph to indicate explicitly when datasets contain outliers.

Histogram: Histograms also emphasize the existence of outliers. Look for isolated bars in the plot.

Scatterplot: We can use Scatterplot to determine outliers in a multivariate setting. Isolated points from regression line gives us the outliers.

3. Using Z-Score:

Z-scores can quantify the unusualness of an observation when your data follow the normal distribution. Z-scores are the number of standard deviations above and below the mean that each value falls. For example, a Z-score of 2 indicates that an observation is two standard deviations above the average while a Z-score of -2 signifies it is two standard deviations below the mean. A Z-score of zero represents a value that equals the mean. Statistical rule of thumb is if Z-Score for a specific data point is less than -4 or greater than 4 is a suspected outlier.

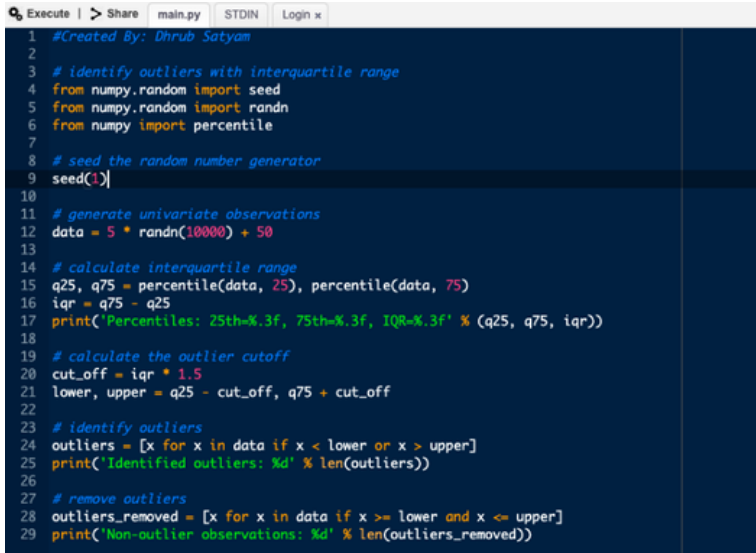
4. Using IQR:

You can use the interquartile range (IQR), several quartile values, and an adjustment factor to calculate boundaries for what constitutes minor and major outliers. Minor and major denote the unusualness of the outlier relative to the overall distribution of values. Major outliers are more extreme. Analysts also refer to these categorizations as mild and extreme outliers.

The IQR is the middle 50% of the dataset. It's the range of values between the third quartile and the first quartile ($Q3 - Q1$). We can take the IQR, $Q1$, and $Q3$ values to calculate the following outlier fences for our dataset:

lower outer, lower inner, upper inner, and upper outer. These fences determine whether data points are outliers and whether they are mild or extreme.

Values that fall inside the two inner fences are not outliers.



```
1 #Created By: Dhrub Satyam
2
3 # identify outliers with interquartile range
4 from numpy.random import seed
5 from numpy.random import randn
6 from numpy import percentile
7
8 # seed the random number generator
9 seed(1)
10
11 # generate univariate observations
12 data = 5 * randn(10000) + 50
13
14 # calculate interquartile range
15 q25, q75 = percentile(data, 25), percentile(data, 75)
16 iqr = q75 - q25
17 print('Percentiles: 25th=%3f, 75th=%3f, IQR=%3f' % (q25, q75, iqr))
18
19 # calculate the outlier cutoff
20 cut_off = iqr * 1.5
21 lower, upper = q25 - cut_off, q75 + cut_off
22
23 # identify outliers
24 outliers = [x for x in data if x < lower or x > upper]
25 print('Identified outliers: %d' % len(outliers))
26
27 # remove outliers
28 outliers_removed = [x for x in data if x >= lower and x <= upper]
29 print('Non-outlier observations: %d' % len(outliers_removed))
```

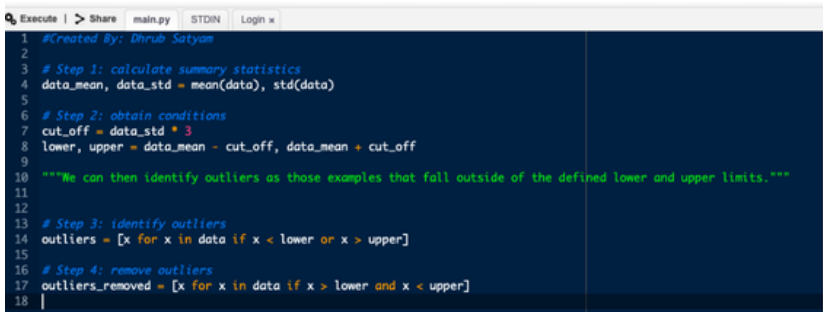
5. Using Hypothesis Testing: You can use hypothesis tests to find outliers. Many outlier tests exist, but I'll focus on one to illustrate how they work. In this post, I demonstrate Grubbs' test, which tests the following hypotheses:

Null: All values in the sample were drawn from a single population that follows the same normal distribution.

Alternative: One value in the sample was not drawn from the same normally distributed population as the other values.

If the p-value for this test is less than your significance level, you can reject the null and conclude that one of the values is an outlier.

6. Standard Deviation Method: We know that if the data follow normal distribution then the data covers 99.7% of the points up to 3 standard deviation. We can have our outliers calculated beyond that on both sides. Python Code for this is as follows:



```
1 #Created By: Dhruv Satyam
2
3 # Step 1: calculate summary statistics
4 data_mean, data_std = mean(data), std(data)
5
6 # Step 2: obtain conditions
7 cut_off = data_std * 3
8 lower, upper = data_mean - cut_off, data_mean + cut_off
9
10 """We can then identify outliers as those examples that fall outside of the defined lower and upper limits."""
11
12
13 # Step 3: identify outliers
14 outliers = [x for x in data if x < lower or x > upper]
15
16 # Step 4: remove outliers
17 outliers_removed = [x for x in data if x > lower and x < upper]
18 |
```

18. When do you reject or accept null hypothesis? List Steps.

Step 1: State the null hypothesis. When you state the null hypothesis, you also have to state the alternate hypothesis. Sometimes it is easier to state the alternate hypothesis first, because that's the researcher's thoughts about the experiment.

Step 2: Support or reject the null hypothesis. Several methods exist, depending on what kind of sample data you have. For example, you can use the P-value method.

If you are able to reject the null hypothesis in Step 2, you can replace it with the alternate hypothesis.

That's it!

When to Reject the Null hypothesis?

Basically, you reject the null hypothesis when your test value falls into the rejection region. There are four main ways you'll compute test values and either support or reject your null hypothesis.

Which method you choose depends mainly on if you have a proportion or a p-value.

(For details refer :

https://www.sagepub.com/sites/default/files/upmbinaries/40007_Chapter8.pdf)

19) What are Type I and Type II errors? How would you proceed to minimise them?

No hypothesis test is 100% certain. Because the test is based on probabilities, there is always a chance of making an incorrect conclusion. When you do a hypothesis test, two types of errors are possible: type I and type II. The risks of these two errors are inversely related and determined by the level of significance and the power for the test. Therefore, you should determine which error has more severe consequences for your situation before you define their risks.

Type I error:

When the null hypothesis is true and you reject it, you make a type I error. The probability of making a type I error is α , which is the level of significance you set for your hypothesis test. An α of 0.05 indicates that you are willing to accept a 5% chance that you are wrong when you reject the null hypothesis. To lower this risk, you must use a lower value for α . However, using a lower value for alpha means that you will be less likely to detect a true difference if one really exists.

Type II error:

When the null hypothesis is false and you fail to reject it, you make a type II error. The probability of making a type II error is β , which depends on the power of the test. You can decrease your risk of committing a type II error by ensuring your test has enough power. You can do this by ensuring your sample size is large enough to detect a practical difference when one truly exists.

The probability of rejecting the null hypothesis when it is false is equal to $1-\beta$. This value is the power of the test.

Decision based on sample	Truth about the population	
	H_0 is true	H_0 is false
Fail to reject H_0	Correct Decision (probability = $1 - \alpha$)	Type II Error - fail to reject H_0 when it is false (probability = β)
Reject H_0	Type I Error - rejecting H_0 when it is true (probability = α)	Correct Decision (probability = $1 - \beta$)

Avoiding Type-I Error:

Minimize the significance level (α) of Hypothesis Test: since we choose the significance level. However, lowering the significance level may lead to a situation wherein the results of the hypothesis test may not capture the true parameter or the true difference of the test.

Avoiding Type-II Error:

- 1.Increase the sample size
- 2.Increase the significance level of Hypothesis Test.

20) What is Bias and Variance? How can we have optimum of both?

Bias: is the difference between the average prediction of our model and the correct value which we are trying to predict. Model with high bias pays very little attention to the training data and oversimplifies the model. It always leads to high error on training and test data.

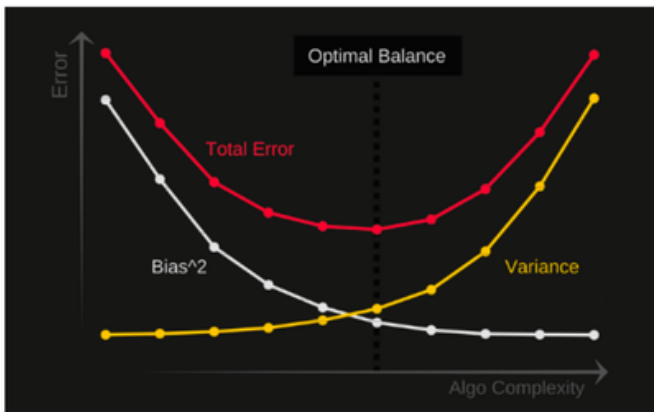
Variance: is the variability of model prediction for a given data point or a value which tells us spread of our data. Model with high variance pays a lot of attention to training data and does not generalize on the data which it hasn't seen before. As a result, such models perform very well on training data but has high error rates on test data.

In supervised learning, underfitting happens when a model unable to capture the underlying pattern of the data. These models usually have high bias and low variance. It happens when we have very less amount of data to build an accurate model or when we try to build a linear model with a nonlinear data.

Also, these kind of models are very simple to capture the complex patterns in data like Linear and logistic regression.

In supervised learning, overfitting happens when our model captures the noise along with the underlying pattern in data. It happens when we train our model a lot over noisy dataset. These models have low bias and high variance. These models are very complex like Decision trees which are prone to overfitting.

Bias Variance Tradeoff:



For any supervised algorithm, having a high bias error usually means it has low variance error and vice versa. To be more specific, parametric or linear ML algorithms often have a high bias but low variance. On the other hand, non-parametric or non-linear algorithms have vice versa.

The goal of any ML model is to obtain a low variance and a low bias state, which is often a task due to the parametrization of machine learning algorithms. Common ways to achieve optimum Bias and Variance are:

1. By minimising total error
2. Using Bagging and resampling techniques
3. Adjusting minor values in Algorithms

(Click to refer : <https://hub.packtpub.com/heres-how-you-can-handle-the-bias-variance-trade-off-in-your-ml-models/>)

21. What is Confusion Matrix?

It is a performance measurement for machine learning classification problem where output can be two or more classes. It is a table with 4 different combinations of predicted and actual values.

		Actual Values	
		Positive (1)	Negative (0)
Predicted Values	Positive (1)	TP	FP
	Negative (0)	FN	TN

Definition of the terms:

Positive (P) : Observation is positive (for example: is an apple).

Negative (N) : Observation is not positive (for example: is not an apple).

True Positive (TP) : Observation is positive, and is predicted to be positive.

False Negative (FN) : Observation is positive, but is predicted negative.

True Negative (TN) : Observation is negative, and is predicted to be negative.

False Positive (FP) : Observation is negative, but is predicted positive.

Classification Rate/Accuracy: Classification Rate or Accuracy is given by the relation:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

However, there are problems with accuracy. It assumes equal costs for both kinds of errors. A 99% accuracy can be excellent, good, mediocre, poor or terrible depending upon the problem.

Recall/Sensitivity/True Positive Rate: Recall can be defined as the ratio of the total number of correctly classified positive examples divide to the total number of positive examples. High Recall indicates the class is correctly recognized (a small number of FN).

$$\text{Recall} = \frac{TP}{TP + FN}$$

Precision: To get the value of precision we divide the total number of correctly classified positive examples by the total number of predicted positive examples. High Precision indicates an example labelled as positive is indeed positive (a small number of FP).

$$\text{Precision} = \frac{TP}{TP + FP}$$

F-measure/F-stats/F1 Score: Since we have two measures (Precision and Recall) it helps to have a measurement that represents both of them. We calculate an F-measure which uses Harmonic Mean in place of Arithmetic Mean as it punishes the extreme values more. The F-Measure will always be nearer to the smaller value of Precision or Recall.

$$F - \text{measure} = \frac{2 * \text{Recall} * \text{Precision}}{\text{Recall} + \text{Precision}}$$

Specificity: Percentage of negative instances out of the total actual negative instances. Therefore, denominator (TN + FP) here is the actual number of negative instances present in the dataset. It is similar to recall but the shift is on the negative instances. Like finding out how many healthy patients were not having cancer and were told they don't have cancer. Kind of a measure to see how separate the classes are.

$$\frac{TN}{TN + FP}$$

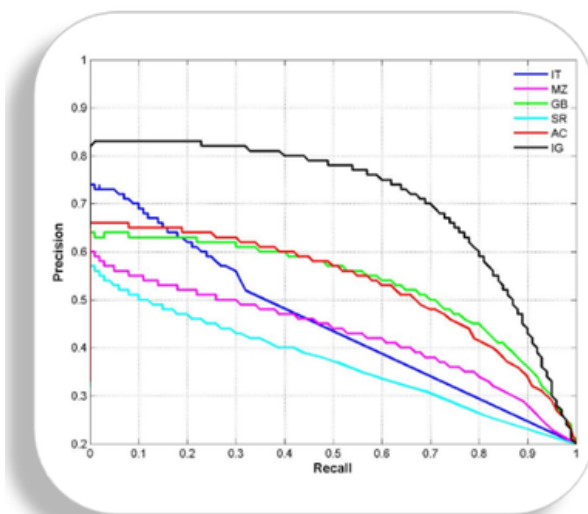
22. What are different ways of measuring the performance of different models?

The different ways are as follows:

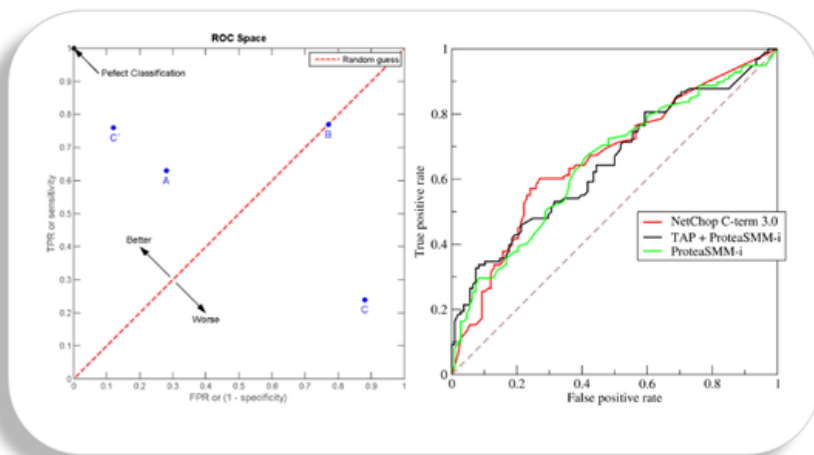
- Confusion matrix
- Accuracy
- Precision
- Recall
- Specificity
- F1 score
- Precision-Recall or PR curve
- ROC (Receiver Operating Characteristics) curve
- PR vs ROC curve

Most of the ways we have already seen in the above question. We will discuss only which are left.

PR Curve: It is the curve between precision and recall for various threshold values. In the figure below we have 6 predictors showing their respective precision-recall curve for various threshold values. The top right part of the graph is the ideal space where we get high precision and recall. Based on our application we can choose the predictor and the threshold value. PR AUC is just the area under the curve. The higher its numerical value the better.



ROC Curve: ROC stands for receiver operating characteristic and the graph is plotted against TPR and FPR for various threshold values. As TPR increases FPR also increases. As you can see in the first figure, we have four categories and we want the threshold value that leads us closer to the top left corner. Comparing different predictors (here 3) on a given dataset also becomes easy as you can see in figure 2, one can choose the threshold according to the application at hand. ROC AUC is just the area under the curve, the higher its numerical value the better.



PR vs ROC Curve:

Both the metrics are widely used to judge a models performance. Which one to use PR or ROC?

The answer lies in TRUE NEGATIVES.

Due to the absence of TN in the precision-recall equation, they are useful in imbalanced classes. In the case of class imbalance when there is a majority of the negative class. The metric doesn't take much into consideration the high number of TRUE NEGATIVES of the negative class which is in majority, giving better resistance to the imbalance. This is important when the detection of the positive class is very important.

Like to detect cancer patients, which has a high class imbalance because very few have it out of all the diagnosed. We certainly don't want to miss on a person having cancer and going undetected (recall) and be sure the detected one is having it (precision).

Due to the consideration of TN or the negative class in the ROC equation, it is useful when both the classes are important to us. Like the detection of cats and dog. The importance of true negatives makes sure that both the classes are given importance, like the output of a CNN model in determining the image is of a cat or a dog.

23. What is correlation? What is the range that a correlation coefficient can have?



Correlation is a statistical technique that can show whether and how strongly pairs of variables are related. For example, height and weight are related; taller people tend to be heavier than shorter people.

The main result of a correlation is called the correlation coefficient (or "r"). It ranges from -1.0 to +1.0. The closer r is to +1 or -1, the more closely the two variables are related.

24. What is ROC Curve? What is AUC?

Refer to Question 16 for ROC Curve.

AUC: ROC curve is a performance measurement for classification problem at various thresholds settings. ROC is a probability curve and AUC represents degree or measure of separability. It tells how much model is capable of distinguishing between classes. Higher the AUC, better the model is at predicting 0s as 0s and 1s as 1s. By analogy, Higher the AUC, better the model is at distinguishing between patients with disease and no disease.

The ROC curve is plotted with TPR against the FPR where TPR is on y-axis and FPR is on the x-axis.

25. What is Collinearity and Correlation?

Correlation measures the strength of linear relationship between two variables.

Collinearity: If in multiple regression analysis, one of the predictors is linearly associated/dependent on other predictor, then this issue is known as collinearity.

For example, let's consider the linear model

$$Y = \alpha x_1 + \beta_1 x_1 + \beta_2 x_2 \dots (1)$$

If predictor x_1 can be expressed as linear combination of x_2 , say, $x_1 = 3 \cdot x_2$

Then this is known as collinearity among the predictors. Note that there will be perfect (or very high) correlation between the predictors as opposed to the assumption of linear regression model (All predictors are assumed to be independent).

Essentially it means that one of the independent variables is not really necessary to the model because its effect/impact on the model is already captured by some of the other variables. This variable is not contributing anything extra to the predictions and can be removed. If we have true collinearity (perfect correlation as in the example above), the one of the predictor is automatically deleted by some of the software's like R, other shows an error or warning for the same.

[Click here for Q.26 to Q.50](#)

To know more information about statistics
and it's implementation in Data Science
check out the links below



IBM Certified Data Science and AI Program for working professionals.

Eligibility: Minimum 1 yrs of exp.

Duration: 7.5 Months

[PDF](#)



IBM Certified Artificial Intelligence and ML Program for working professionals.

Eligibility: Minimum 4 yrs of exp.

Duration: 9 Months

[PDF](#)



IBM Certified Data Science and AI Program for Managers and Leaders.

Eligibility: Minimum 8 yrs of exp.

Duration: 11 Months

[PDF](#)

Talk to our Industrial Expert and start your journey
towards Data Science and AI

[Schedule](#)